

VidToMe: Video Token Merging for Zero-Shot Video Editing

Xirui Li¹ Chao Ma¹ Xiaokang Yang¹ Ming-Hsuan Yang²
¹Shanghai Jiao Tong University ²UC Merced

Project webpage: <https://vidtome-diffusion.github.io>

Abstract

Diffusion models have made significant advances in generating high-quality images, but their application to video generation has remained challenging due to the complexity of temporal motion. Zero-shot video editing offers a solution by utilizing pre-trained image diffusion models to translate source videos into new ones. Nevertheless, existing methods struggle to maintain strict temporal consistency and efficient memory consumption. In this work, we propose a novel approach to enhance temporal consistency in generated videos by merging self-attention tokens across frames. By aligning and compressing temporally redundant tokens across frames, our method improves temporal coherence and reduces memory consumption in self-attention computations. The merging strategy matches and aligns tokens according to the temporal correspondence between frames, facilitating natural temporal consistency in generated video frames. To manage the complexity of video processing, we divide videos into chunks and develop intra-chunk local token merging and inter-chunk global token merging, ensuring both short-term video continuity and long-term content consistency. Our video editing approach seamlessly extends the advancements in image editing to video editing, rendering favorable results in temporal consistency over state-of-the-art methods.

1. Introduction

Diffusion models [8, 9, 16, 38, 46, 47] have made significant advances in synthesizing media content, allowing for the creation of diverse, high-quality images. However, diffusion models have yet to achieve high quality in generating videos. Due to the complexity of temporal motion in videos, training a video diffusion model requires a massive amount of data and computation resources. To avoid learning temporal motion from scratch, zero-shot video editing leverages a pre-trained image diffusion model to translate a source video into a new one, retaining motion from the source video. Separately editing each frame likely results in inconsistent frames (Fig 1 Per-frame Editing). Ex-

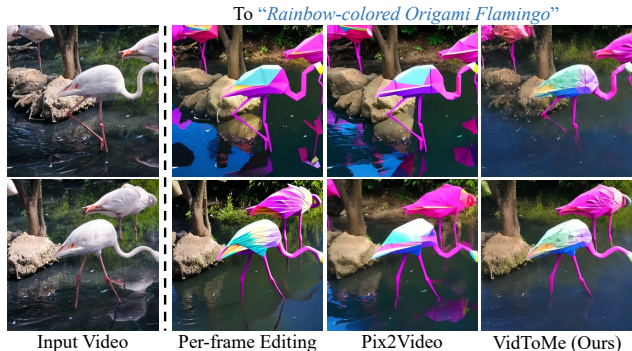


Figure 1. Given an input source video and a text prompt, we leverage a pre-trained image diffusion model [38] to edit the video. The state-of-the-art zero-shot video editing approaches, e.g., Pix2Video [6], struggle to generate temporal consistent frames with self-attention extension. Our proposed method merges tokens across frames, rendering higher temporal consistency.

isting video editing methods [6, 34, 51, 52, 54] typically extend the self-attention modules of diffusion models to process multiple frames jointly instead of separately. Despite the promise, two issues ensue with such approaches. First, though cross-frame attention encourages a roughly consistent appearance, the generated frames lack strict consistency in details. As human perception is sensitive to video continuity, tiny changes or jittering between frames can significantly degrade the quality of generated videos. Second, including multi-frame tokens in self-attention increases memory consumption quadratically. Computing self-attention on four frames requires 16 times larger GPU memory than on one frame. With these limitations, state-of-the-art video editing methods such as Pix2Video [6] struggle to generate temporal consistent videos, as shown in Fig. 1. Thus, it is imperative to develop effective and efficient diffusion-based zero-shot video editing methods.

In this work, we present a novel method, VidToMe, to enhance the temporal consistency of generated videos by merging the tokens in diffusion models across video frames. Our motivation comes from the recent developments [4, 5] in compressing transformer tokens to improve computational efficiency. For video editing, we observe that the



Figure 2. Comparison of frames edited to “Pop Art Style” by PnP [48] with or without VidToMe. Left: Edit results. Right: Visualized token matching between two frames as flow maps. Color represents the direction of the matched token in another frame. We label the denoising timestamp above ($1000 \rightarrow 0$). Our method aligns correspondent tokens and fixes the inconsistencies (window, clothes) in per-frame editing.

transformer tokens of video frames are much more correlated in the temporal domain than in the spatial domain. We can align tokens across frames according to the correlation and compress the temporally redundant tokens to facilitate joint self-attention. We show that when multiple frames share the same set of tokens in the self-attention module, the diffusion model naturally generates temporally consistent video frames. Hence, we propose merging tokens over time to compress and unify the internal diffusion feature embeddings, achieving temporal consistency in generated videos and reducing memory consumption in computing self-attention across frames.

Our method involves merging tokens in one frame with the most similar ones in another frame. As shown in Fig. 2, this merging strategy allows us to match and align tokens according to the temporal correspondence between frames. Thus, applying VidToMe fixes the misalignment of details in per-frame editing. As processing all frames at once is not feasible, we divide the video into chunks and perform intra-chunk local token merging and inter-chunk global token merging, ensuring both short-term and long-term consistency. Short-term consistency improves video continuity, while long-term consistency prevents the video content from drifting over time. Note that our video editing method can be seamlessly integrated with existing controlling schemes [48, 56], taking full advantage of the advancements in image editing for video editing. Extensive experiments show that the proposed video editing method performs well regarding temporal consistency and text alignment over the state-of-the-art approaches. The main contributions of this work are three-fold:

- We propose a novel diffusion-based approach, VidToMe, to merge self-attention tokens across frames when generating video frames, improving temporal

consistency and computational efficiency.

- We design a video editing pipeline that jointly generates all video frames with short-term local token merging and long-term global token merging to enforce feature alignment throughout the video.
- We comprehensively evaluate our method to show the state-of-the-art video editing performance.

2. Related Work

As several thorough surveys on image and video generation [1, 7, 8, 25] exist in the literature, here we discuss diffusion models and related image and video editing schemes.

Diffusion-based Image and Video Synthesis. Diffusion Models (DM) [16, 45, 47] have recently achieved state-of-the-art performance in numerous tasks, including image generation [8, 9, 30, 38, 46]. DMs learn to reverse a forward diffusion process and generate an image by gradually denoising it from pure noise. Notable examples of improving DMs include [15, 20, 43] and numerous applications [2, 11, 21, 24, 26–28, 40]. Benefiting from large-scale pretraining [35, 44], text-to-image DMs have shown impressive results in generating high-quality and diverse images [31, 36, 38, 42]. Naturally, DMs have been applied to video synthesis, typically by incorporating temporal layers into image DMs [3, 17, 18]. Despite the demonstrated success in unconditional video generation [18, 55], text-to-video DMs are not as satisfying as image ones. Due to the complexity of temporal motion, training video DMs requires intensive computation resources and large-scale annotated video datasets, which significantly hinders the progress of this field.

Diffusion-based Image Editing. In addition to text, some works have introduced additional control signals for image editing or controllable image generation [41, 48, 56]. Some schemes introduce adapter layers [29] or other trainable modules [50, 56] to accept additional control signals. ControlNet [56] supports various conditions such as edge maps, depth maps, and key points by finetuning an attached copy of DM. Other methods edit a source image by manipulating intermediate diffusion features [14, 48] or optimization-based guidance [10, 32]. Plug-and-Play [48] maintains image structure by injecting self-attention maps and internal features from the source image. Self-guidance [10] and pix2pix-zero [32] edit the image by imposing a guidance loss optimized during generation. StableDiffusion2 [38] presents a depth-conditioned model that directly includes the depth map in its input. In this paper, we perform video editing by applying these image editing methods to video frames while preserving temporal coherence via merging video tokens.

Diffusion-based Video Editing. With the recent success of text-to-image DMs in powering text-driven image editing [14, 27, 48], many works apply a pre-trained text-to-

image DM [38] for text-driven video editing. The critical problem is how to keep temporal coherency in generation. Tune-A-Video [52] inflates the DM with temporal attention layers and finetunes on the source video. vid2vid-zero [51] maintains the video structure by injecting cross-attention maps from the source video. In [6], Pix2Video guides the generation with a reference frame by self-attention features injection and latent update. Rerender-A-Video [54] fuses the previous frame warped by the source video optical flow and applies multi-stage latent operations. On the other hand, Fate/Zero [34] uses a dedicated attention blending block to inject attention maps from the source video. TokenFlow [13] shares a similar idea with our method to enforce temporal consistency by unifying self-attention tokens. It computes the inter-frame correspondences by extracting tokens from the source video. Then the tokens are propagated between the jointly-edited keyframes according to the correspondances. Note that these methods commonly extend the self-attention modules into the spatial-temporal domain to encourage consistent appearance across frames. However, extending self-attention does not enforce temporal consistency well and increases memory overhead. Our method simultaneously addresses these two problems by merging similar tokens across video frames.

3. Preliminaries

Latent Diffusion Model. Diffusion models [16, 45, 47] are a class of generative models based on an iterative denoising process. An image DM supposes a forward process where a clean image x_0 is corrupted by Gaussian noise ϵ ,

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon, \quad (1)$$

where $t = 1, \dots, T$ is the current timestep and $\{\alpha_t\}$ are the monotonically decreasing noise schedule. Then, starting from random Gaussian noise, DM reverses the forward process to generate an image by estimating the noise direction and progressively denoising it.

Recent large-scale diffusion models [36, 38, 42] operate in the latent space to improve performance and efficiency. These latent diffusion models train an autoencoder [23] to map the image between pixel and latent space. Let $\mathcal{E}(\cdot)$ and $\mathcal{D}(\cdot)$ be the encoder and the decoder, where $\mathcal{E}(x) = z, \mathcal{D}(z) \approx x$. Both the training and inference are conducted in the latent space. Typically, a UNet [39] ϵ_θ is trained to estimate the noise with the objective

$$\min_{\theta} E_{z, \epsilon \sim \mathcal{N}(0, I), t} \|\epsilon - \epsilon_\theta(z_t, t, c)\|, \quad (2)$$

where c is the text embedding in text-to-image DMs. In this work, we base our experiments on Stable Diffusion [38], a large-scale text-to-image latent diffusion model.

Token Merging. Token Merging (ToMe) [5] is a method to increase the throughput of existing ViT models by gradually merging redundant tokens in the transformer blocks.

It combines similar tokens to reduce the redundancy as well as the number of tokens, speeding up the computation. Our method leverages the lightweight bipartite soft matching algorithm of ToMe to merge tokens across video frames.

Given input tokens T , the algorithm first partitions the tokens into a source (src) and destination (dst) set and computes the pair-wise cosine similarity between the two sets. Then, src tokens are linked to their most similar token in dst . Next, r most similar edges are selected, and connected tokens are merged. Finally, all the tokens are concatenated as the output. We use the dst token as the merged token value instead of averaging the value of merged tokens, which produces better results in practice. Our method divides merged tokens after self-attention to keep the token number unchanged. Like [4], we divide a token simply by assigning its value to all restored tokens, which means a token merged from two tokens will be separated into two identical tokens. We define the token merging and unmerging operations, M and U :

$$\begin{aligned} E &= \text{Match}(src, dst, r), \\ T_m &= M(T, E), T_u = U(T_m, E). \end{aligned} \quad (3)$$

$\text{Match}(\cdot)$ outputs the matching map E with r edges from src to dst . $M(\cdot)$ and $U(\cdot)$ merge and unmerge tokens according to matching E .

4. Proposed Method

Our objective is to generate an edited video that matches a given editing prompt while preserving the motion and structure of a source video. To achieve this, we use a pre-trained text-to-image diffusion model to generate individual frames. We apply DDIM inversion [46] and existing controlling methods [38, 48, 56] for image editing to preserve the source frame structure. However, more effort is required to achieve temporal consistency across video frames. We observe that transformer tokens are correlated across frames as the temporal correspondence in videos. Thus, we compress multi-frame tokens by merging similar tokens together so that the self-attention module extracts consistent features for each frame. The unified internal features promote the diffusion model to generate consistent video frames. Fig. 3 presents an overview of the proposed method.

The proposed video token merging strategy, VidToMe, is detailed at the bottom of Fig. 3. We first merge tokens across frames in one video chunk to enforce short-term video continuity. The locally merged tokens are combined with a set of global tokens from previous chunks, enabling long-term token sharing. Joint self-attention extracts consistent features on merged tokens, which are then propagated to each frame by token unmerging. Our video token merging algorithm has two advantages. First, merged

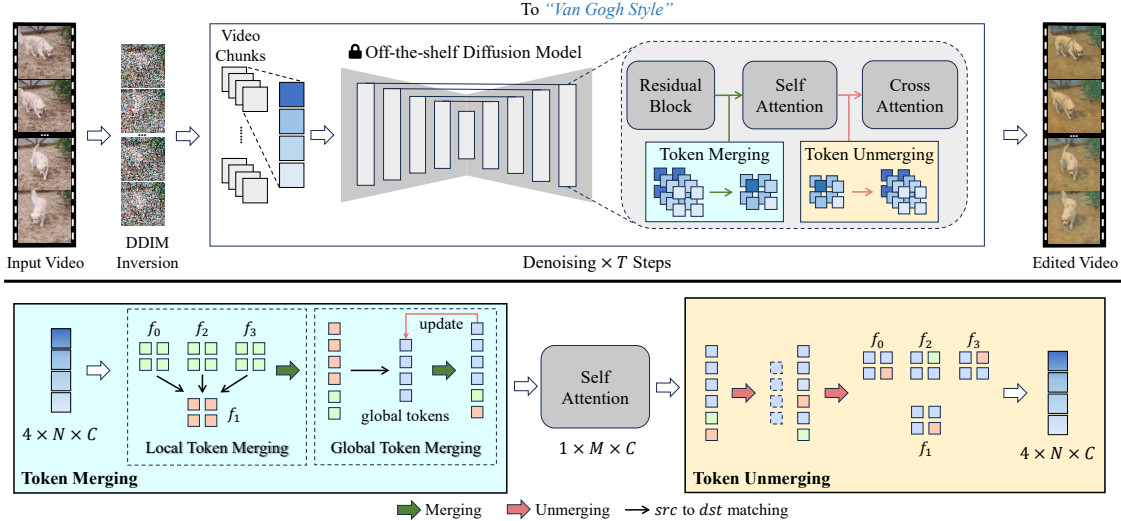


Figure 3. Pipeline of our proposed video editing method, VidToMe. We apply DDIM inversion on the source video frames to obtain the initial noisy latents. We denoise frame latents with an off-the-shelf text-to-image diffusion model, combining an existing controlling method [48, 56]. In each iteration, frames are split into chunks and denoised by the diffusion model, where we attach our lightweight token merging and unmerging operation around the self-attention modules. We first merge tokens locally into a random frame in the chunk. Then, the merged tokens are combined with the global tokens maintained across chunks in one iteration. After self-attention, we unmerge the tokens to the original size for the following operations.

tokens are shared across frames, enforcing temporal consistency. Second, token merging compresses redundant tokens in the self-attention module, improving computation efficiency. As a lightweight parameter-free operation, VidToMe can seamlessly integrate with existing image editing methods [38, 48, 56]. We first provide an overview of our video editing pipeline in Section 4.1 and then elaborate on our video token merging method in Section 4.2.

4.1. Video Editing Pipeline

Given the source video \mathcal{V} with n frames (f_1, f_2, \dots, f_n) , we first invert each frame to noise by DDIM inversion. DDIM inversion applies the deterministic DDIM sampler [46] in the reversed order $t = 0 \rightarrow T$, turning a clean latent z_0 into a noisy latent z_T . Using the inverted frames as the initial noise, we iteratively denoise them by an off-the-shelf text-to-image diffusion model with edit prompt \mathcal{P} as the text condition.

Unlike existing methods [6, 54] that generate each frame separately, we process all frames together in each denoising iteration. At iteration t , we randomly split frames $(z_t^1, z_t^2, \dots, z_t^n)$ into a sequence of chunks (C_1, C_2, \dots, C_m) . Each chunk contains B consecutive frames except for the first and last chunks. The initial chunk $C_1 = (z_t^1, \dots, z_t^b)$ where b is a random integer in $[1, B]$. The randomness here ensures the probability of frames split into a chunk proportional to the time interval, avoiding the “chunks” seen in the generated video. We process each

chunk of frames by the diffusion model in a random order.

In the diffusion model, frames are concatenated in the batch dimension, which means the model treats all frames as separate images. To enforce temporal consistency, we merge tokens across frames in the self-attention module of the diffusion model. A diffusion model is typically realized as a UNet [39], composed of a series of downsample and upsample blocks. A layer in each block consists of a residual block, a self-attention module [49], and a cross-attention module. Among them, the self-attention module has been shown to be highly correlated to the structure and appearance of the image [48]. Thus, we merge multi-frame tokens before self-attention so that the self-attention jointly processes the merged tokens and outputs consistent features. Note that our method only changes the input tokens of the self-attention without any modification to the self-attention operation. To perform the following processes, we restore the output tokens to the original size by token unmerging. As merging tokens in deep blocks may degrade the generation quality, we perform token merging in the first two downsample blocks and the last two upsample blocks.

After T iterations of denoising, we obtain the edited frames latents $(z_0^1, z_0^2, \dots, z_0^n)$ and the edited video \mathcal{V}^* after VAE decoding. Our method works with an existing controlling method for image editing to preserve the structure of source frames, such as ControlNet [56], Plug-and-Play (PnP) [48], and depth-conditioned diffusion model [38].

4.2. Video Token Merging

This section presents our video token merging strategy, which focuses on the self-attention module in the diffusion model. A self-attention module takes a sequence of input tokens and outputs the same number of tokens. The input and output tokens are denoted as T_{in} and T_{out} , both belonging to the space $R^{B \times N \times C}$, where B is the number of frames, N is the number of tokens per frame, and C is the feature dimension. We enforce temporal coherence by merging input tokens across B frames in one video chunk (local token merging) and merging global tokens from previous chunks of frames (global token merging). After self-attention, we unmerge the output tokens to their original size. These operations are performed on the input and output tokens without modifying the self-attention module.

Local Token Merging. Given a set of input tokens denoted by $T_{in} = \{T_{in}^f\}_{f=0}^{B-1}$, we randomly select one out of B frames as the current target frame, *e.g.*, the k^{th} frame. We then apply the bipartite soft matching algorithm mentioned in Section 3 and Equation 3 to merge the other frames to the target frame:

$$T_{lm} = M(T_{in}, \text{Match}(T_{in}^{src}, T_{in}^{dst}, r)),$$

where $T_{in}^{src} = \{T_{in}^f\}_{f=0, f \neq k}^{B-1}$ and $T_{in}^{dst} = T_{in}^k$. We set $r = p(B-1)N$ where $(B-1)N$ is the *src* token number and p is the merging ratio. A large merging ratio (*e.g.*, $p = 0.9$) can be used as video frames are highly redundant. Local token merging enforces consistency in a small frame chunk.

However, for long-term consistency, we need more than short-term video continuity. For example, a video’s first and last frames will never be processed in one chunk, leading to appearance drifting along the video. Thus, we need another way to share tokens across the whole video. Enlarging the chunk size or implementing a hierarchical merging is helpful but requires an even larger memory capacity. Instead, we propose a simple yet effective global token merging strategy for long-term consistency.

Global Token Merging. At each iteration, we maintain a set of global tokens denoted as T_g that spans across video chunks. The initial global tokens are set to be the locally merged tokens of the first chunk, *i.e.*, $T_g^1 = T_{lm}^1$. For the k^{th} frame chunk, we merge its locally merged tokens T_{lm}^k with the previous global tokens T_g^{k-1} as the following operation:

$$T_{gm} = M(\{T_{lm}^k, T_g^{k-1}\}, \text{Match}(T_{lm}^k, T_g^{k-1}, r)), \quad (4)$$

where T_{gm} represents the final input to the self-attention module. In practice, we randomly assign *src* and *dst* to local and global tokens. We can update the global tokens to include the tokens from the current frames in several ways. One possible option is to use merged tokens T_{gm} as new global tokens. However, this approach is not feasible for arbitrarily long videos as it always increases the number of

global tokens. Instead, we unmerge T_{gm} back to local tokens T_{lm}^u and global tokens and set the current global tokens to be the unmerged local tokens, *i.e.*, $T_g^k = T_{lm}^u$.

Self-Attention Analysis. We analyze the self-attention operation on merged tokens in more detail. The input, denoted as T_{gm} , comprises M tokens. We can infer $M = (0.11B + 0.99)N = 1.43N$, assuming chunk size $B = 4$ and merging ratio $p = 0.9$ for local and global merging. The tokens are mapped to Q, K, V matrices during self-attention. For multiplication between Q and K , the original input of size $4 \times N \times C$ has a space complexity of $O(4N^2)$. In contrast, the complexity is reduced to half with the merged tokens $T_{gm} \in R^{1 \times M \times C}$ as input, $O(M^2) \approx O(2N^2)$. Essentially, our token merging method combines multiple frames into one, reducing redundancy among frames. Self-attention then identifies consistent features in this unified frame.

Token Unmerging. The output tokens T_o of the self-attention module need to be restored to their original shape as separate frames to perform the following image operations. As such, we first unmerge the tokens into local and global tokens and then unmerge the local tokens into B separate frames, reversing the merging process. Denoting respective matching maps for local and global token merging as E_l and E_g , we formulate the token unmerging as $U(T_o, E_g) = (T_{local}, T_{global})$ to divide the local tokens and $U(T_{local}, E_l) = T_{out}$ to obtain the final output. The unmerged tokens in the output are identical to the original merged tokens, ensuring consistency across frames.

5. Experimental Results

5.1. Experiment Setting

Our method performs video editing with a pretrained text-to-image model and an existing image-controlling method. In this work, we use Stable Diffusion (SD) [38] (version 1.5) as the image generator, and DDIM scheduler [46] with sampling step $T = 50$ for inversion and sampling. For the controlling method, we combine our method with Plug-and-Play (PnP) [48], ControlNet [56], and SD2-Depth. The parameters are chunk length $B = 4$ and merging ratio $p = 0.9$ and 0.8 for local and global merging. More results and videos are available in the supplementary material. All the source code and datasets will be released.

Dataset. Similar to prior works [6, 34, 53], we select 20 videos from DAVIS [33] as source videos for performance evaluation. These videos include a range of motion from slow to fast and feature various subjects such as humans, vehicles, and animals. We edit each of the 20 videos with three types of prompts: (i) Style prompts edit the global style. (ii) Object prompts edit the object’s appearance and attributes. (iii) Background prompts change the video background. To obtain edit prompts, we use some prompts from [53] and

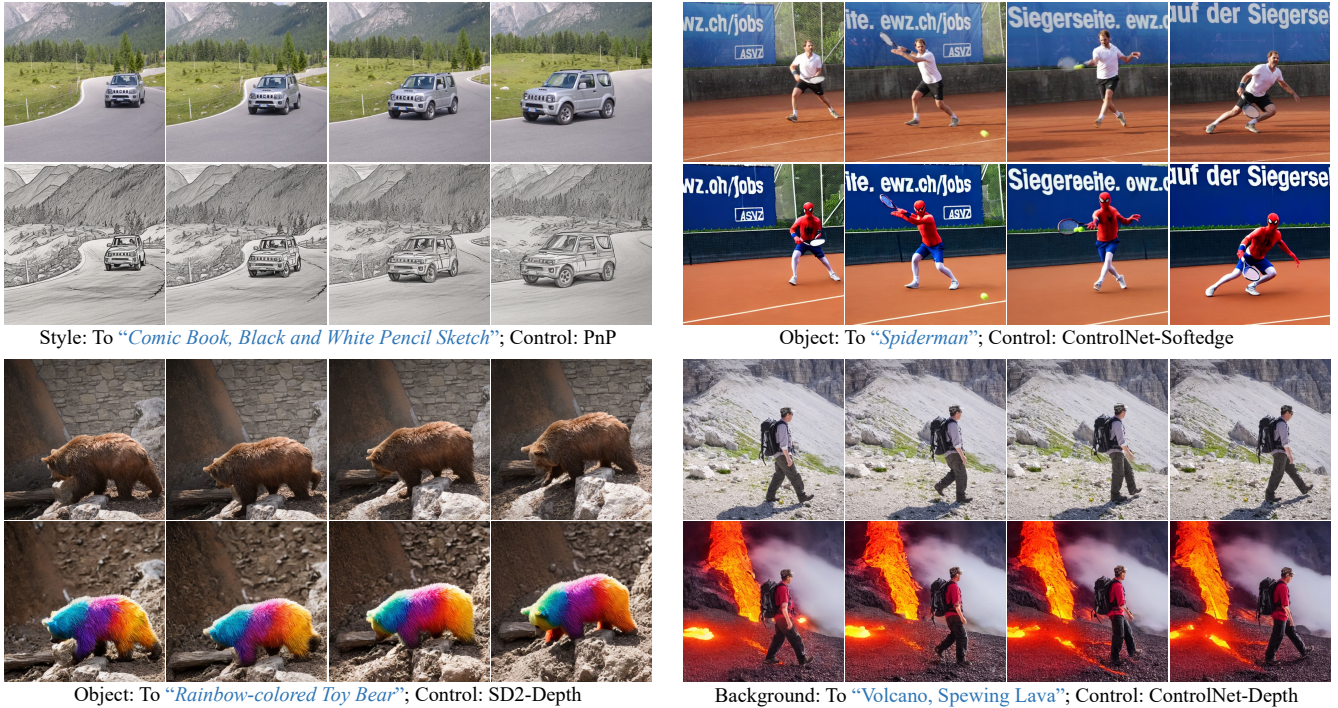


Figure 4. Sample editing results of our method. Our method can seamlessly integrate with existing controlling methods [38, 48, 56] to temporal consistently edit videos in various aspects. We label each sample with the edit prompt and the applied controlling method.

generate the others using GPT-3.5 [37]. We generate 60 edited videos for evaluation, each containing 32 frames with a resolution of 512×512 .

Metrics. We note that existing video editing methods use different metrics to evaluate the editing performance. In our experiment, we incorporate metrics used by prior works [6, 34] and new measures. We assess the editing performance based on three key criteria: (i) Temporal Consistency: This includes Interpolation Error and PSNR [19], Warp Error, and Frame CLIP Score. (ii) Text Alignment: This includes Directional CLIP Score [12] and Text CLIP Score. (iii) User Study. Based on the Interpolation Error and PSNR used to evaluate video interpolation performance in previous studies [19], we propose to measure the video continuity by interpolating a target frame by its previous and next frames and computing the root-mean-squared (RMS) difference and PSNR between the estimated and target frames. The metric better reflects the generated video continuity itself without relying on the source video optical flow. We use the Directional CLIP Score [12] to evaluate the consistency between image and prompt change. It computes the cosine similarity in CLIP space between the difference in frames and the difference in prompts from source to edit. For more details about the metrics, please refer to the supplementary material. We conduct a user study with 10 out of 60 edited videos. Users choose their preferred video from the results edited by both baselines and our method.

The user preference rate is used as the final metric.

Baselines. We evaluate our method against four state-of-the-art video editing techniques: Text2Video-Zero [22], Tune-A-Video [52], vid2vid-zero [51], and Pix2Video [6]. All the methods are implemented using default settings except for vid2vid-zero. We have to turn off its Spatial-Temporal attention that includes all frames in the self-attention, which is infeasible to fit the GPU memory (40GB) when processing 32-frame videos. Text2Video-Zero allows zero-shot text-to-video generation, and we apply it with depth control to perform video editing. Tune-A-Video fine-tunes the model on the source video frames before sampling the edited video. We use StableDiffusion v1.5 for the first three methods and StableDiffusion2-Depth [38] for Pix2Video as it requires a depth-conditioned model by default. It is worth noting that all the baseline methods use some self-attention extension that includes multiple frames.

5.2. Main Results

Qualitative Evaluation. Fig. 5 compares our editing results with baseline methods on evaluation videos. While Text2Video-Zero [22] produces high-quality frames, it lacks continuity between them. The edited frames by Text2Video-Zero do not align with the source frame in appearance due to the random initial noise it used. Tune-A-Video [52] struggles to learn the motion of the source video and fails when the motion is complex. The edited frames

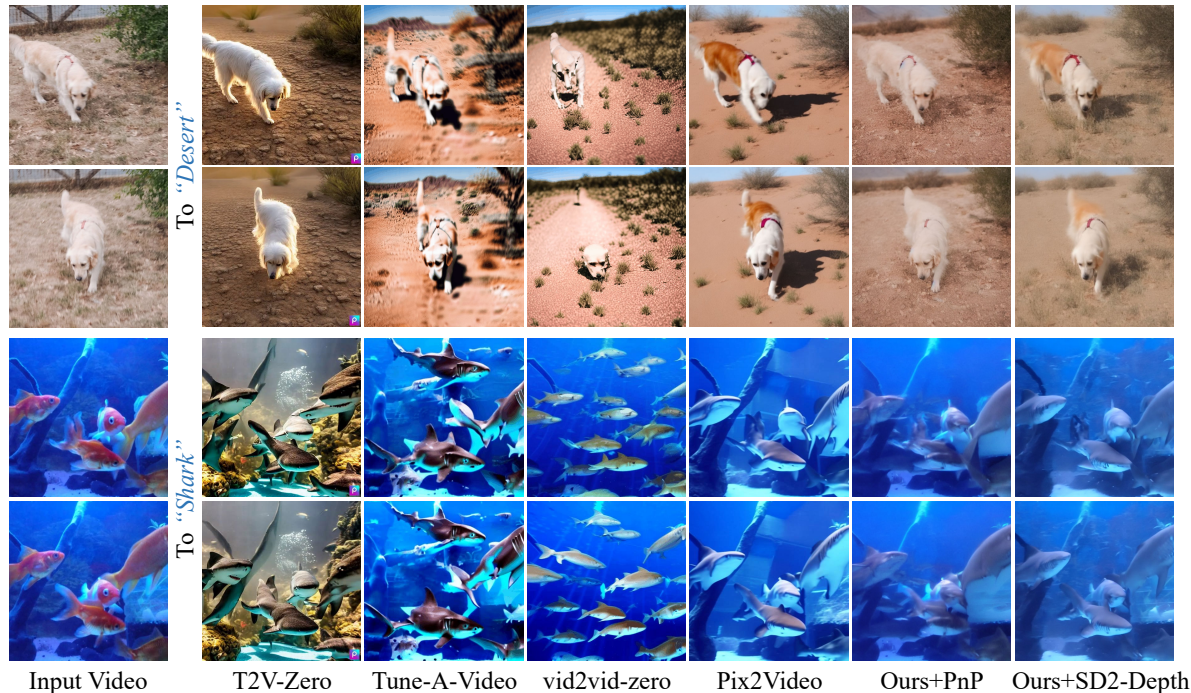


Figure 5. Qualitative comparison of our method and baselines. The editing results of our method are consistent over time in global style and local texture and preserve the source frame structure well.

	Temporal Consistency				Text Alignment		User Study
	Inp. Err.↓	Inp. PSNR↑	Warp Err.↓	Frame C.s.↑	Directional C.s.↑	Text C.s.↑	Preference Rate↑
Text2Video-Zero [22]	0.203	19.4	0.094	0.969	0.139	0.282	0.167
Tune-A-Video [52]	0.235	18.1	0.068	0.957	0.132	0.273	0.089
vid2vid-zero [51]	0.219	18.6	0.049	0.957	0.118	0.270	0.030
Pix2Video* [6]	0.139	22.5	0.020	0.968	0.166	0.285	0.170
Ours + ControlNet-Depth [56]	0.154	22.1	0.026	0.973	0.159	0.284	
Ours + PnP [48]	0.111	25.0	0.013	0.975	0.140	0.271	0.544
Ours + SD2-Depth [38]*	0.105	25.6	0.012	0.971	0.168	0.282	-

Table 1. Quantitative evaluation results. Red and blue indicates the best and second-best result. *: Use the same base model as Pix2Video, SD2-Depth [38]. Others use SDv1.5. C.s.: CLIP Score. Inp. Err.: Interpolation Error.

by Tune-A-Video contain wave-like jittering, degrading the video quality. vid2vid-zero [51] generates unstable video and fails to preserve the frame structure. Pix2Video [6] achieves good consistency between edited frames but generates unnatural jittering and blurring results. In contrast, VidToMe generates consistent frames that adhere to the edit prompt while preserving the source frame structure. Fig. 4 showcases more sample editing results. Our method can integrate seamlessly with existing controlling methods, providing users with more control over the editing process and enabling video editing in various aspects.

Quantitative Evaluation. We present quantitative evaluation results in Table 1. The first three baselines do not perform well in terms of temporal consistency. Text2Video-Zero [22] uses random noise instead of noise inverted from

source frames, which results in frames that are different from the source frames, as shown in the qualitative study. Though some subjects appreciate the diversity in its results, its continuity is unsatisfactory. Tune-A-Video [52] and vid2vid-zero [51] are not preferred by subjects due to temporal inconsistency. Pix2Video [6] achieves higher temporal consistency than the other baselines. It uses SD2 [38], which has better generation fidelity than SDv1.5 and thus gets a high directional and text CLIP Score. However, its results still suffer from jittering and lack of long-term consistency. Our proposed VidToMe achieves better temporal consistency and text alignment than the first three baselines. It also outperforms Pix2Video regarding temporal consistency when using the same base model, SD2-Depth. Furthermore, Our editing results are preferred by over half of

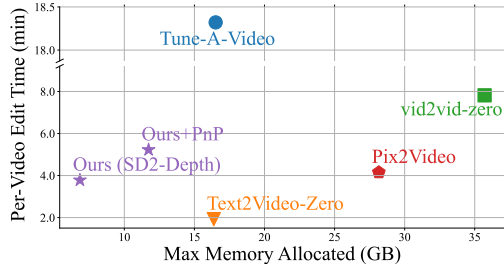


Figure 6. Editing efficiency comparison. We consider both time (time to perform one editing) and space (max GPU memory allocated in editing) efficiency. Test on one NVIDIA Tesla A100.

	Inp. Err. ↓	Directional C.s. ↑
Per-frame S.A.	0.253	0.165
Extended S.A.	0.140	0.168
VidToMe S.A. (Ours)	0.105	0.168

Table 2. Performance comparison of different multi-frame self-attention (S.A.) operations.

the subjects. Our method offers flexibility in balancing consistency and alignment by using different control methods such as ControlNet-Depth [56] or PnP [48].

5.3. Ablation Studies

Efficiency Analysis. In Fig. 6, we compare the editing efficiency of different methods with their default setting. Most other methods require either high memory capacity or a long editing time. Although Text2Video-Zero [22] edits videos quickly without noise inversion, its performance is poor. On the other hand, our proposed method, VidToMe, reduces memory consumption with video token merging while generating videos quickly. With a minimal memory consumption of less than 7 GB, our method can run on some personal devices.

Multi-frame Self-Attention. We ablate on multi-frame self-attention choices in video editing. In per-frame self-attention, each frame is processed separately, leading to inconsistencies. On the other hand, extending self-attention to include multi-frame tokens is adopted by most existing methods. This method enables cross-frame attention and produces better consistency. Our approach merges tokens to enforce temporal consistency, achieving the best performance without sacrificing the editing effect.

Token Merging Operation. We ablate on our token merging choice in Fig. 7. The original ToMe algorithm merges tokens by averaging their values. However, we find that this can lead to a lack of diversity and randomness in the generated videos, such as the flamingo features being single-colored. Therefore, we directly replace the value of the merged tokens with *dst* tokens. Global token merging is crucial for keeping long-term consistency in videos. With-

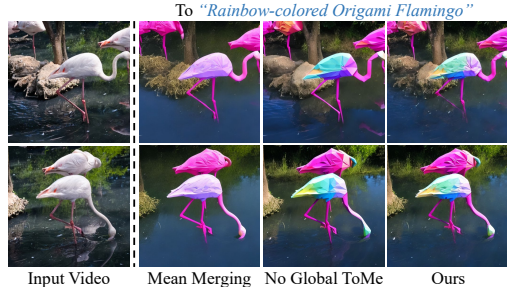


Figure 7. Ablation on Token Merging Operations. Merging tokens by mean instead of replacement reduces the edit fidelity. Without global token merging, the feather color changes.

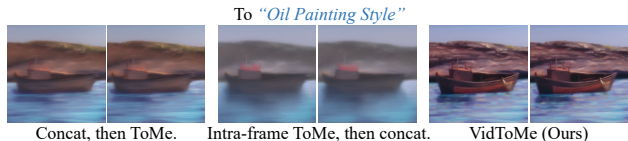


Figure 8. Ablation on local token merging strategy. Try: (i) Concatenate multi-frame tokens and then perform vanilla ToMe. (ii) Perform ToMe in each frame, then concatenate their tokens. (iii) Our local token merging strategy. We control the token number after merging to be the same.

out it, our method fails to maintain consistent rainbow feather colors throughout the video.

Local Token Merging Strategy. There are several possible strategies to merge tokens through multiple consecutive frames. One straightforward approach is to perform the original ToMe [4] among tokens from all frames, where *dst* tokens are randomly selected. Or we can apply ToMe in each frame first and then concatenate all merged tokens. Fig. 8 shows that these two methods result in blurry frames, while our local merging strategy preserves the quality of the generated frames.

Limitations. Our method has two main limitations. First, the editing capability of our method depends on the performance of the selected image editing technique. If the editing technique fails on a single frame, our method also fails to edit the entire video. Second, although our similarity-based matching performs well in most cases, it has room for improvement. Objects with similar features are sometimes incorrectly merged and mixed in the output results. We plan to explore a more precise token-matching approach in the future to address these issues.

6. Conclusion

This work proposes a diffusion-based zero-shot video editing method, VidToMe. Our approach unifies and compresses internal diffusion features by matching and merging tokens across video frames in the self-attention module during generation, resulting in temporally consistent edited

video frames. We implement VidToMe as lightweight token merging and unmerging blocks attached to the self-attention module, making it compatible with any existing image editing method and diffusion models.

References

- [1] Nuha Aldausari, Arcot Sowmya, Nadine Marcus, and Gelareh Mohammadi. Video generative adversarial networks: a review. *ACM Computing Surveys (CSUR)*, 55(2): 1–25, 2022. 2
- [2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *CVPR*, pages 18208–18218, 2022. 2
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *CVPR*, pages 22563–22575, 2023. 2
- [4] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. *CVPR Workshop on Efficient Deep Learning for Computer Vision*, 2023. 1, 3, 8
- [5] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your ViT but faster. In *ICLR*, 2023. 1, 3
- [6] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *ICCV*, pages 23206–23217, 2023. 1, 3, 4, 5, 6, 7
- [7] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A Bharath. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1):53–65, 2018. 2
- [8] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *IEEE TPAMI*, 2023. 1, 2
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 34:8780–8794, 2021. 1, 2
- [10] Dave Epstein, Allan Jabri, Ben Poole, Alexei A Efros, and Aleksander Holynski. Diffusion self-guidance for controllable image generation. *arXiv preprint arXiv:2306.00986*, 2023. 2
- [11] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2022. 2
- [12] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. 6
- [13] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 3
- [14] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *ICLR*, 2022. 2
- [15] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 2
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. 1, 2, 3
- [17] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- [18] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022. 2
- [19] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation. In *CVPR*, pages 9000–9008, 2018. 6
- [20] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *NeurIPS*, 35:26565–26577, 2022. 2
- [21] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *NeurIPS*, 35: 23593–23606, 2022. 2
- [22] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. 6, 7, 8
- [23] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *ICLR*, 2014. 3
- [24] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022. 2
- [25] Ming-Yu Liu, Xun Huang, Jiahui Yu, Ting-Chun Wang, and Arun Mallya. Generative adversarial networks for image and video synthesis: Algorithms and applications. *Proceedings of the IEEE*, 109(5):839–862, 2021. 2
- [26] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*, pages 11461–11471, 2022. 2
- [27] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jianjun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2021. 2
- [28] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *CVPR*, pages 6038–6047, 2023. 2
- [29] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhonggang Qi, Ying Shan, and Xiao-hu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2
- [30] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, pages 8162–8171. PMLR, 2021. 2

- [31] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, pages 16784–16804. PMLR, 2022. 2
- [32] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. In *SIGGRAPH*, pages 1–11, 2023. 2
- [33] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016. 5
- [34] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv:2303.09535*, 2023. 1, 3, 5, 6
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 2
- [36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 3
- [37] Partha Pratim Ray. Chatgpt: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems*, 2023. 6
- [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1, 2, 3, 4, 5, 6, 7
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015. 3, 4
- [40] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, pages 22500–22510, 2023. 2
- [41] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *SIGGRAPH*, pages 1–10, 2022. 2
- [42] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35:36479–36494, 2022. 2, 3
- [43] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022. 2
- [44] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 35:25278–25294, 2022. 2
- [45] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265. PMLR, 2015. 2, 3
- [46] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2020. 1, 2, 3, 4, 5
- [47] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICML*, 2020. 1, 2, 3
- [48] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *CVPR*, pages 1921–1930, 2023. 2, 3, 4, 5, 6, 7, 8
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017. 4
- [50] Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. Sketch-guided text-to-image diffusion models. In *SIGGRAPH*, pages 1–11, 2023. 2
- [51] Wen Wang, kangyang Xie, Zide Liu, Hao Chen, Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-shot video editing using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023. 1, 3, 6, 7
- [52] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *ICCV*, pages 7623–7633, 2023. 1, 3, 6, 7
- [53] Jay Zhangjie Wu, Xiuyu Li, Difei Gao, Zhen Dong, Jinbin Bai, Aishani Singh, Xiaoyu Xiang, Youzeng Li, Zuwei Huang, Yuanxi Sun, Rui He, Feng Hu, Junhua Hu, Hai Huang, Hanyu Zhu, Xu Cheng, Jie Tang, Mike Zheng Shou, Kurt Keutzer, and Forrest Iandola. Cvpv 2023 text guided video editing competition. *arXiv preprint arXiv:2310.16003*, 2023. 5
- [54] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. In *SIGGRAPH*, 2023. 1, 3, 4
- [55] Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in projected latent space. In *CVPR*, pages 18456–18466, 2023. 2
- [56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 2, 3, 4, 5, 6, 7, 8